# Quantized Distributed Gradient Tracking Algorithm with Linear Convergence in Directed Networks

Yongyang Xiong, Ligang Wu, *Fellow, IEEE*, Keyou You, *Senior Member, IEEE*, and Lihua Xie, *Fellow, IEEE*

*Abstract*—Communication efficiency is a major bottleneck in the applications of distributed networks. To address the problem, the problem of quantized distributed optimization has attracted a lot of attention. However, most of the existing quantized distributed optimization algorithms can only converge sublinearly. To achieve linear convergence, this paper proposes a novel quantized distributed gradient tracking algorithm (Q-DGT) to minimize a finite sum of local objective functions over directed networks. Moreover, we explicitly derive lower bounds for the number of quantization levels, and prove that Q-DGT can converge linearly even when the exchanged variables are respectively quantized with 3 quantization levels. Numerical results also confirm the efficiency of the proposed algorithm.

*Index Terms*—Quantized communication, distributed optimization, gradient tracking algorithm, directed networks.

## I. INTRODUCTION

RECENT years have witnessed tremendous progress in distributed optimization due to its wide applications in formation control [1], distributed resource allocation [2], online optimization [3], localization systems [4], game theory [5], to name a few. They require a group of networked nodes to cooperatively optimize the sum of their local cost functions via local communications. A comprehensive review of this topic can be found in [6], [7].

Although distributed algorithms are capable of solving complex tasks in a collaborative manner, limited communication capacity is a major bottleneck in distributed networks, especially for large-scale distributed machine learning. How to design communication-efficient distributed algorithms has attracted an increasing attention [8]–[10]. For instance, the encoding-decoding scheme in [11] has been designed to distributedly solve linear equations [12], distributed optimization problems [13], [14]. To further reduce the size of data transmission, the recent work [15] showed that the sign of relative state between neighbors is sufficient for achieving convergence. As errors are inevitable for a finite-precision quantizer, the QDGD algorithm proposed in [16] achieves vanishing consensus error even in the presence of non-vanishing noise by modifying the contribution of the received quantized information for each node. By incorporating quantization scheme into the push-sum algorithm [17], the authors of [18] proposed distributed algorithms over directed networks for both convex and non-convex functions. Since the aforementioned works are derived by the distributed gradient descent (DGD) [17],

Y. Xiong and K. You are with the Department of Automation, Beijing National Research Center for Information Science and Technology, Tsinghua University, Beijing 100084, P. R. China. E-mail: xiongyy@tsinghua.edu.cn; youky@tsinghua.edu.cn.

L. Wu is with the Department of Control Science and Engineering, Harbin Institute of Technology, Harbin 150001, P. R. China. E-mail: ligangwu@hit.edu.cn.

L. Xie is with the School of Electrical and Electronic Engineering, Nanyang Technological University, 50 Nanyang Avenue, Singapore 639798. E-mail: elhxie@ntu.edu.sg.

TABLE I: DISTRIBUTED OPTIMIZATION ALGORITHMS

| References | digraphs | linear convergence | 1-bit communication |
|---|---|---|---|
| [9], [13], [16] | ✗ | ✗ | ✗ |
| [8], [14], [15] | ✗ | ✗ | ✓ |
| [19]–[22], [31], [34] | ✗ | ✓ | ✗ |
| [18], [23], [24] | ✓ | ✗ | ✗ |
| [2], [25]–[27] | ✓ | ✓ | ✗ |
| Our work | ✓ | ✓ | ✓ |

they can only achieve sublinear convergence even for the strongly convex functions.

How to accelerate the convergence speed is fundamentally important to reduce communication cost. Recently, a few significant efforts have been devoted to designing quantized distributed algorithms with linear convergence. For instance, ref. [19], [20] proposed DQOA and LEAD, respectively, under the assumption that the randomized quantizer $Q(\cdot)$ is an unbiased and $\delta$-contracted operator, i.e. $\mathbb{E}[Q(x)] = x$ and $\mathbb{E}\|Q(x) - x\|^2 \leq \delta\|x\|^2$ for all $x \in \mathbb{R}^m$. Clearly, this assumption excludes some important quantizers, e.g., the binary quantizer. Ref. [21] proposed Q-NEXT by dynamically adjusting the center of the quantization interval. Ref. [22] established a trade-off between the convergence speed and the communication cost per iteration so that linear convergence can be guaranteed. Although the aforementioned quantized algorithms [19]–[22] converge linearly, they are designated only for undirected networks. Note that extending distributed algorithms from undirected networks to directed networks is non-trivial [23]–[27]. In fact, if the directed network is unbalanced, i.e., there exists at least a node that the sum of the weights of its outgoing nodes is not equal to that of its incoming nodes (see e.g., [28], [29]), the DGD finally minimizes a weighted average of local functions. Hence, an additional variable is usually exchanged between nodes to eliminate the effects of the unbalancedness [24], [27]. To resolve the unbalancedness issue, the push-pull/$\mathcal{AB}$ algorithm [25], [26] and its variant [2] leverage row-stochastic matrix and column-stochastic matrix simultaneously and achieve exact linear convergence for strongly convex and smooth functions. In sharp contrast to the subgradient-based quantized algorithms in [13], [14], [18] that only the decision variable needs to be quantized, the quantizer cannot be directly incorporated into push-pull/$\mathcal{AB}$ as it will result in an accumulation of quantization errors [30], thereby the convergence cannot be guaranteed.

A question naturally arises: Whether it is possible to develop a quantized distributed algorithm over directed networks that converges linearly even for one-bit communication? In this paper, we give a positive answer. A comparison of our work with the state-of-the-art works is provided in Table I. The main contributions of this work are summarized as follows:

1) We propose a novel quantized distributed algorithm Q-DGT over directed networks. The Q-DGT is remarkably robust to quantization errors, and achieves linear convergence.

2) We explicitly provide the lower bounds of the quantization levels to resolve the saturation issue for the finite-level quantizers,

which even supports the extreme 3-level quantization.

The remainder of this paper is organized as follows. We formulate the problem in Section II. The Q-DGT is provided in Section III. Section IV includes the convergence analysis. Simulation results are presented in Section V. In Section VI, we conclude this paper.

**Notation.** We use $x_i$ to denote the $i$-th element of vector $x$; $\mathbf{1}_n(\mathbf{0}_n)$ denotes a column vector with its all elements equaling to one(zero). The notation $f = \mathcal{O}(h)$ means there exists a positive constant $\upsilon < \infty$ such that $f \leq \upsilon h$. $\nabla F(x(k)) \triangleq (\nabla f_1^{\mathrm{T}}(x_1(k)), ..., \nabla f_n^{\mathrm{T}}(x_n(k)))^{\mathrm{T}}$. For an arbitrary vector norm $\| \cdot \|$, the induced norm of a matrix $W = (w_1, ..., w_m) \in \mathbb{R}^{n \times m}$ is defined as $\|W\| = \sqrt{\sum_{i=1}^{m} \|w_i\|^2}$. Throughout, we slightly abuse the notation of vector norms and their induced matrix norms for simplicity.

## II. PRELIMINARIES AND PROBLEM FORMULATION

In this section, we first introduce some basics of graph theory. In what follows we formulate the problem of interest.

### A. Basics of Graph Theory

Consider a digraph $\mathcal{G}(\mathcal{V}, \mathcal{E})$, where $\mathcal{V} = \{1, ..., n\}$ denotes the set of nodes, $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$ represents the set of directed links, and $(i, j) \in \mathcal{E}$ implies that node $j$ can receive information from node $i$. We denote $\mathcal{N}_i^{\mathrm{in}} = \{j : (j, i) \in \mathcal{E}\} \cup \{i\}$ and $\mathcal{N}_i^{\mathrm{out}} = \{j : (i, j) \in \mathcal{E}\} \cup \{i\}$ as the in-neighbor set and out-neighbor set of node $i$, respectively. A digraph is strongly connected if there exists a directed path between any pair of distinct nodes, which is commonly used in the literature [6].

*Assumption 1:* The digraph $\mathcal{G}(\mathcal{V}, \mathcal{E})$ is strongly connected.

### B. Problem Formulation

Consider the digraph $\mathcal{G}(\mathcal{V}, \mathcal{E})$ where each node $i \in \mathcal{V}$ privately processes a convex function $f_i : \mathbb{R}^m \to \mathbb{R}$. All nodes collaboratively solve the following optimization problem[1]:

$$\underset{x \in \mathbb{R}^m}{\text{minimize}} \ f(x) = \sum_{i=1}^{n} f_i(x). \tag{1}$$

In such a problem, each node $i$ maintains a local estimate $x_i(k) \in \mathbb{R}^m$ of the decision vector $x$ at each step $k$ and can only share its own information with a subset of nodes via the communication network. We make the following assumptions on the local functions:

*Assumption 2:* Each $f_i$ is $\mu$-strongly convex, i.e., there exists a $\mu > 0$, such that $f_i(y) \geq f_i(x) + \nabla f_i(x)^{\mathrm{T}}(y - x) + \frac{\mu}{2}\|x - y\|_2^2$, $\forall x, y \in \mathbb{R}^m$.

*Assumption 3:* Each $f_i$ is $L$-smooth, i.e., $\|\nabla f_i(x) - \nabla f_i(y)\|_2 \leq L\|x - y\|_2$ for some $L > 0$, $\forall x, y \in \mathbb{R}^m$.

Assumptions 2-3 are standard for the linear convergence in literature, see e.g. [31], [34]. Under Assumption 2, the problem (1) has a unique optimal solution $x^\star \in \mathbb{R}^m$.

The main objective of this paper is to design a distributed algorithm where nodes are only allowed to communicate quantized variables over $\mathcal{G}(\mathcal{V}, \mathcal{E})$, with linear convergence to the exact optimal solution $x^\star$ of problem (1).

## III. ALGORITHM DEVELOPMENT

In this section, we first explain why a quantizer cannot be directly incorporated into the push-pull/$\mathcal{AB}$ Algorithm [25], [26]. Then, we propose the Q-DGT and show that it is robust to quantization errors. Finally, we introduce the quantization rule.

[1]For clarity of presentation, we only consider the scalar variable case, i.e., $m = 1$, as our algorithm and its analysis can be easily extended to the vector case by using the Kronecker operator.

### A. Push-pull/$\mathcal{AB}$ Algorithm with Naive Quantization Does Not Work

In the push-pull/$\mathcal{AB}$ algorithm [25], [26], each node $i \in \mathcal{V}$ maintains two vectors $x_i(k)$ and $y_i(k)$ per step $k \in \mathbb{N}$, and performs the following updates:

$$x_i(k+1) = \sum_{j=1}^{n} a_{ij}x_j(k) - \eta y_i(k), \tag{2a}$$

$$y_i(k+1) = \sum_{j=1}^{n} b_{ij}y_j(k) + \nabla f_i(x_i(k+1) - \nabla f_i(x_i(k)), \tag{2b}$$

where $y_i(0) = \nabla f_i(x_i(0))$, $\mathcal{A} = [a_{ij}]_{n \times n}$ and $\mathcal{B} = [b_{ij}]_{n \times n}$ are weight matrices induced by $\mathcal{G}(\mathcal{V}, \mathcal{E})$ satisfying: (1) $a_{ij} > 0$ for $j \in \mathcal{N}_i^{\mathrm{in}}$, otherwise $a_{ij} = 0$, and $\sum_{j \in \mathcal{N}_i^{\mathrm{in}}} a_{ij} = 1$; (2) $b_{ij} > 0$ for $i \in \mathcal{N}_j^{\mathrm{out}}$, otherwise $b_{ij} = 0$, and $\sum_{i \in \mathcal{N}_j^{\mathrm{out}}} b_{ij} = 1$. Naive quantization means that $x_j(k)$ and $y_j(k)$ in (2) are replaced by their quantized versions $\hat{x}_j(k)$ and $\hat{y}_j(k)$, respectively. That is, if $x_j(k)$ and $y_j(k)$ in (2) are directly quantized as $\hat{x}_j(k)$ and $\hat{y}_j(k)$, respectively. Then,

$$x_i(k+1) = \sum_{j=1}^{n} a_{ij}\hat{x}_j(k) - \eta y_i(k), \tag{3a}$$

$$y_i(k+1) = \sum_{j=1}^{n} b_{ij}\hat{y}_j(k) + \nabla f_i(x_i(k+1) - \nabla f_i(x_i(k)), \tag{3b}$$

where $\sigma_{x_j}(k) \triangleq \hat{x}_j(k) - x_j(k)$ and $\sigma_{y_j}(k) \triangleq \hat{y}_j(k) - y_j(k)$ are the quantization errors. Note that if $\mathcal{A} = \mathcal{B}$, then (3) is exactly the quantized algorithm in [32]. However, taking summation over $i \in \mathcal{V}$, (3b) implies that

$$\mathbf{1}_n^{\mathrm{T}}y(k+1) = \mathbf{1}_n^{\mathrm{T}}\nabla F(x(k+1)) + \sum_{l=0}^{k} \mathbf{1}_n^{\mathrm{T}}\sigma_y(l). \tag{4}$$

Thus the quantization errors are accumulated, i.e., $\sum_{l=0}^{k} \mathbf{1}_n^{\mathrm{T}}\sigma_y(l)$. No matter whether $\sigma_y(k)$ in (4) converges or not, $\mathbf{1}_n^{\mathrm{T}}y(k)$ cannot exactly track the global gradient $\mathbf{1}_n^{\mathrm{T}}\nabla F(x(k))$.

This observation was first pointed out in [30] and then the author proposed a robust push-pull algorithm. However, the work [30] does not involve the design of quantizer and simply assumes that $\mathbb{E}[\sigma_x(k)] = \mathbb{E}[\sigma_y(k)] = \mathbf{0}_n$, and $\mathbb{E}[\|\sigma_x(k)\|^2] \leq \sigma_x$, $\mathbb{E}[\|\sigma_y(k)\|^2] \leq \sigma_y$ for some $\sigma_x, \sigma_y > 0$. This condition is clearly not satisfied for the deterministic quantizers. In addition, the algorithm in [30] can only converge to a neighborhood of the optimal solution in expectation. All above motivates us to propose the Q-DGT.

### B. The Q-DGT Algorithm

In this work, we design a dynamic encoding-decoding scheme for quantized communication (see Fig. 1). At step $k$, each node $j \in \mathcal{V}$ encodes $x_j(k)$ into $r_j(k)$ by using:

$$r_j(k) = Q_{K_x}\left(\frac{1}{h(k)}(x_j(k) - \hat{x}_j(k-1))\right), \tag{5}$$

where $\hat{x}_j(-1) = \mathbf{0}_m$ and $\hat{x}_j(k-1)$ is an estimation of $x_j(k-1)$, $h(k)$ is a decaying scaling function. Note that we quantize the scaled "innovation", i.e., $\frac{1}{h(k)}(x_j(k) - \hat{x}_j(k-1))$. The reason is that the amplitude of the prediction error is usually smaller than that of the state itself such that the scaled "innovation" can be quantized by fewer bits. However, it brings challenge for the finite-level quantizer to avoid saturation. We will show later that the value of $x_j(k) - \hat{x}_j(k-1)$ decays to zero at the speed of the same order of $h(k)$, and rigorously prove that the scaled "innovation" can always be upper bounded by a finite constant. Then, node $j$ broadcasts $r_j(k)$ to its out-neighbors. Upon $r_j(k)$ is received by the out-neighbor node $i \in \mathcal{N}_j^{\mathrm{out}}$, it decodes $r_j(k)$ as follows:

$$\hat{x}_j(k) = h(k)r_j(k) + \hat{x}_j(k-1). \tag{6}$$

Fig. 1: The encoding-decoding scheme.

Here $h(k)$ plays a critical role in estimating the states of node $j$. We highlight that all the out-neighbors of node $j$ receive the same information, so we do not distinguish the specific subscript. The above encoding-decoding scheme is performed for $y_j(k)$ in the same way, i.e., encode $y_j(k)$ into $s_j(k)$ and decode $s_j(k)$ to $\hat{y}_j(k)$. However, the deterministic quantization errors makes it infeasible to apply the robust push-pull algorithm [30] directly in our setting (see Section III-A). To resolve it, we design the updates of node $i \in \mathcal{V}$ as follows:

$$x_i(k+1) = x_i(k) + \alpha \sum_{j=1}^n a_{ij}\left(\hat{x}_j(k) - \hat{x}_i(k)\right)$$
$$- \eta\left(y_i(k) - y_i(k-1)\right), \tag{7a}$$

$$y_i(k+1) = (1-\beta)y_i(k) + \beta \sum_{j=1}^n b_{ij}\hat{y}_j(k)$$
$$+ \nabla f_i(x_i(k+1)), \tag{7b}$$

where $\alpha, \beta \in (0,1)$ are two positive constants, $\eta \geq 0$ is a constant step size that will be specified later. Although node $i \in \mathcal{V}$ can access to its true values $x_i(k)$ and $y_i(k)$ at step $k$, the estimate $\hat{x}_i(k)$ and $\hat{y}_i(k)$ are also used in our algorithm for error compensations, which is of the similar spirit as in the quantized average consensus in [11]. We summarize the Q-DGT in Algorithm 1.

Now, we demonstrate why the Q-DGT is robust to quantization errors. Let $\mathcal{A}_\alpha \triangleq (1-\alpha)I_n + \alpha\mathcal{A}$ and $\mathcal{B}_\beta \triangleq (1-\beta)I_n + \beta\mathcal{B}$. Then, (7) can be rewritten as the following compact form:

$$x(k+1) = \mathcal{A}_\alpha x(k) + \alpha(\mathcal{A} - I_n)\sigma_x(k)$$
$$- \eta(y(k) - y(k-1)), \tag{8a}$$

$$y(k+1) = \mathcal{B}_\beta y(k) + \nabla F(x(k+1)) + \epsilon_y(k), \tag{8b}$$

where $\epsilon_{y_i}(k) \triangleq \beta \sum_{j=1}^n b_{ij}\sigma_{y_j}(k)$.

Let $z(k) \triangleq y(k) - y(k-1)$. Then,

$$x(k+1) = \mathcal{A}_\alpha x(k) + \alpha(\mathcal{A} - I_n)\sigma_x(k) - \eta z(k), \tag{9a}$$
$$z(k+1) = \mathcal{B}_\beta z(k) + (\nabla F(x(k+1)) + \epsilon_y(k))$$
$$- (\nabla F(x(k)) + \epsilon_y(k-1)). \tag{9b}$$

Assumption 1 implies that $\mathcal{A}$ has a unique nonnegative left eigenvector $\pi_\mathcal{A}$ such that $\pi_\mathcal{A}^T \mathbf{1}_n = 1$ and $\pi_\mathcal{A}^T \mathcal{A} = \pi_\mathcal{A}^T$, and $\mathcal{B}$ has a unique nonnegative right eigenvector $\pi_\mathcal{B}$ such that $\pi_\mathcal{B}^T \mathbf{1}_n = 1$ and $\mathcal{B}\pi_\mathcal{B} = \pi_\mathcal{B}$ [2]. Define $\bar{x}(k) \triangleq \pi_\mathcal{A}^T x(k)$ and $\bar{z}(k) \triangleq \mathbf{1}_n^T z(k)$, we obtain

$$\bar{x}(k+1) = \bar{x}(k) - \eta\pi_\mathcal{A}^T z(k), \tag{10a}$$
$$\bar{z}(k+1) = \bar{z}(k) + \mathbf{1}_n^T\left(\nabla F(x(k+1)) + \epsilon_y(k)\right)$$
$$- \mathbf{1}_n^T\left(\nabla F(x(k)) + \epsilon_y(k-1)\right). \tag{10b}$$

Conducting mathematical induction for (10b) yields that

$$\bar{z}(k+1) = \mathbf{1}_n^T\left(\nabla F(x(k+1)) + \epsilon_y(k)\right). \tag{11}$$

Notably, the accumulated error $\sum_{l=0}^k \mathbf{1}_n^T \sigma_y(l)$ in (4) disappears in (11). If $\mathbf{1}_n^T \epsilon_y(k)$ tends to zero, then $\bar{z}(k+1)$ tends to the exact global gradient $\mathbf{1}_n^T \nabla F(x(k+1))$. In contrast to [30], we do not make any

---

**Algorithm 1** The Q-DGT —from the view of node $i$

1: **Initialization:** randomly initialize $x_{i,0}$, and $y_{i,0}$ for each $i \in \mathcal{V}$.
2: **for** $k = 0, 1, 2, ...$ **do**
3:     **Encoder:** calculate $r_i(k)$ and $s_i(k)$.
4:     **Communication:** broadcast $r_i(k)$ and $s_i(k)$ to its out-neighbors, and receive $r_j(k)$ and $s_j(k)$ from its in-neighbors $j \in \mathcal{N}_i^{\text{in}}$.
5:     **Decoder:** calculate $\hat{x}_j(k)$ and $\hat{y}_j(k)$.
6:     **Updation:** update $x_i(k+1)$ and $y_i(k+1)$ via (7).
7: **end for**
8: **Return**: $\{x_i(k)\}$.

---

assumption on the error $\epsilon_y(k)$. This requires to design the Q-DGT (7) carefully and handle the joint effects of quantization errors on $x(k)$ and $z(k)$. Specifically, our algorithm can converge linearly and even support 3-level quantization.

### C. The Quantization Rule

The uniform quantizer $Q_K(\cdot)$ for a vector $u = (u_1, ..., u_m)^T$ is defined as $Q_K(u) = (q(u_1), ..., q(u_m))^T$ with

$$q(u_i) = \begin{cases} 0, & -1/2 < u_i \leq 1/2 \\ k, & \frac{2k-1}{2} < u_i \leq \frac{2k+1}{2}, \ k = 1, ..., K \\ K, & u_i > \frac{2K+1}{2} \\ -q(-u_i), & u_i \leqslant -1/2 \end{cases}$$

for $i = 1, ..., m$. The quantizer $q(\cdot)$ maps a real number to a finite set $\mathcal{S} = \{0, \pm k; k = 1, 2, ..., K\}$ with $K \in \mathbb{N}_+$. The quantization level of $q(\cdot)$ is $2K+1$. If $\|u\|_\infty \leq K+1/2$, the quantizer is not saturated, and the quantization error is bounded, i.e., $\|u - Q_K(u)\|_\infty \leqslant 1/2$.

## IV. CONVERGENCE ANALYSIS

In this section, we first establish lower bounds for the quantization levels to solve the saturation issue. Then, the linear convergence of Q-DGT under finite-level quantization is rigorously proved. Finally, we show that Q-DGT converges linearly even with 3-level quantization.

### A. Design of Finite Quantization Levels to Avoid Saturation

Note that the joint effect of quantization on the evolutions of $x(k)$ and $z(k)$ brings challenges to design the finite quantization levels. To solve this issue, we first derive the upper bound of the feasible step size, and then obtain the lower bounds for the quantization levels.

*Lemma 1 ( [2], [25]):* Suppose Assumption 1 holds. There exists matrix norms $\|\cdot\|_\mathcal{A}$ and $\|\cdot\|_\mathcal{B}$ such that $\sigma_\mathcal{A} \triangleq \|\mathcal{A}_\alpha - \mathbf{1}_n\pi_\mathcal{A}^T\|_\mathcal{A} < 1$ and $\sigma_\mathcal{B} \triangleq \|\mathcal{B}_\beta - \pi_\mathcal{B}\mathbf{1}_n^T\|_\mathcal{B} < 1$. Moreover, there exists positive scalars $\delta_{\mathcal{A}2}$, $\delta_{\mathcal{B}2}$, $\delta_{\mathcal{A}\mathcal{B}}$ and $\delta_{\mathcal{B}\mathcal{A}}$ such that for any $X \in \mathbb{R}^{n \times p}$, we have $\delta_{\mathcal{B}\mathcal{A}}^{-1}\|X\|_\mathcal{B} \leq \|X\|_\mathcal{A} \leq \delta_{\mathcal{A}\mathcal{B}}\|X\|_\mathcal{B}$, $\|X\|_2 \leq \|X\|_\mathcal{B} \leq \delta_{\mathcal{B}2}\|X\|_2$, and $\delta_{\mathcal{A}2}^{-1}\|X\|_\mathcal{A} \leq \|X\|_2 \leq \|X\|_\mathcal{A}$.

Define

$$\Theta(k) \triangleq \left(\|\bar{x}(k) - x^\star\|_2, \|x(k) - \mathbf{1}_n\bar{x}(k)\|_\mathcal{A}, \|z(k) - \pi_\mathcal{B}\bar{z}(k)\|_\mathcal{B}\right)^T.$$
$$\tag{12}$$

To facilitate the subsequent analysis, we further define: $\kappa_1 \triangleq \|I_n - \mathbf{1}_n\pi_{\mathcal{A}}^{\mathrm{T}}\|_{\mathcal{A}}$, $\kappa_2 \triangleq \|\pi_{\mathcal{B}}\|_{\mathcal{A}}$, $\kappa_3 \triangleq \|I_n - \pi_{\mathcal{B}}\mathbf{1}_n^{\mathrm{T}}\|_{\mathcal{B}}$, $\kappa_4 \triangleq \|\mathcal{A}_\alpha - I_n\|_2$. The following lemma provides a linear matrix inequality, which will be instrumental in establishing the lower bound for quantization level.

*Lemma 2:* Suppose Assumptions 1-3 hold. If the step size $\eta \leq \frac{1}{(\mu+L)\pi_{\mathcal{A}}^{\mathrm{T}}\pi_{\mathcal{B}}}$, then

$$\Theta(k+1) \preceq G\Theta(k) + \varsigma(k), \tag{13}$$

where $\Theta(k)$ is defined in (12), the notation $\preceq$ means element-wise less than or equal to, $G \in \mathbb{R}^{3\times3}$ and $\varsigma(k) \in \mathbb{R}^{3\times1}$ are given by (25) and (26), respectively.

*Proof:* See Appendix A. ∎

In Lemma 2, the presence of $\varsigma(k)$ is due to quantization errors. If $\varsigma(k)$ linearly converges to $\mathbf{0}_3$, then we can prove that $\Theta(k)$ linearly converges to $\mathbf{0}_3$ provided that the spectral radius $\rho(G) < 1$. After that, the linear convergence of Q-DGT can be proved. We first provide a sufficient condition in terms of the step size $\eta$ to guarantee $\rho(G) < 1$.

*Lemma 3:* Suppose Assumptions 1-3 hold. If the step size $\eta$ satisfies

$$\eta \leq \min\left\{\frac{1}{(\mu+L)\pi_{\mathcal{A}}^{\mathrm{T}}\pi_{\mathcal{B}}}, \frac{1-\sigma_{\mathcal{A}}}{2\sqrt{n}\kappa_1\kappa_2 L\delta_{\mathcal{A}2}}, \frac{1-\sigma_{\mathcal{B}}}{2\delta_{\mathcal{B}2}\kappa_3 L}, \frac{2\Gamma_3}{\Gamma_2 + \sqrt{\Gamma_2^2 + 4\Gamma_1\Gamma_3}}\right\}, \tag{14}$$

where $\Gamma_i$, $i = 1, 2, 3$, are constants given by (28). Then, $\rho(G) < 1$.

*Proof:* See Appendix B. ∎

*Remark 1:* Note that the network information is required to calculate the upper bound of step size. If $\mathcal{A}$ and $\mathcal{B}$ are known, then the parameters $\kappa_1$-$\kappa_4$, the step size $\eta$, and the spectral radius $\rho(G)$ can be obtained. Typically, this requirement is necessary even for unquantized push-pull/$\mathcal{A}\mathcal{B}$ algorithms [25], [26], [30].

It is known that if $\rho(G) < 1$, then there exist a matrix norm $\|\cdot\|_{\mathcal{C}}$ and a constant $\tau$ such that $\|G^k\|_2 \leq \tau\hat{\rho}^k$ for an arbitrarily small constant $\varpi > 0$ [33], where

$$\hat{\rho} \triangleq \rho(G) + \varpi < 1. \tag{15}$$

Now, we are in a position to provide conditions on the quantization levels, under which the saturation issue can be solved.

*Theorem 1:* Suppose Assumptions 1-3 hold. Let $h(k) = C\xi^k$, where $C$ is a known positive constant, and $\xi \in (\hat{\rho}, 1)$. The step size $\eta$ is chosen according to (14). Then the quantizers will never saturate provided that $K_x$ and $K_y$ satisfy the following conditions:

$$K_x \geq \max\left\{\frac{v_1}{C} - \frac{1}{2}, \frac{\sqrt{3}\varphi_1\|\Theta(0)\|_2}{C\xi} + \frac{2\alpha n+1}{2\xi} - \frac{1}{2}, \right.$$
$$\left. \frac{\sqrt{3}\varphi_1\tau\|\Theta(0)\|_2}{C\xi}\bar{\Upsilon} + \frac{2\alpha n+1}{2\xi} + \frac{n\eta\beta}{2\xi^2} - \frac{1}{2}\right\},$$

$$K_y \geq \max\left\{\frac{v_2}{C} - \frac{1}{2}, \frac{\sqrt{3}\varphi_2\tau\|\Theta(0)\|_2}{C}\bar{\Upsilon} + \frac{n\beta+1}{2\xi} - \frac{1}{2}\right\}, \tag{16}$$

where $v_1 \triangleq \max_{i\in\mathcal{V}}\|x_i(0)\|_\infty$, $v_2 \triangleq \max_{i\in\mathcal{V}}\|y_i(0)\|_\infty$, $\varphi_1$ and $\varphi_2$ are given in (33), $\bar{\Upsilon} \triangleq 1 + \frac{\tilde{\varsigma}\hat{\rho}}{\xi(\xi-\hat{\rho})\|\Theta(0)\|_2} + \frac{\tilde{\varsigma}}{\xi\tau\|\Theta(0)\|_2}$ with the constant $\tilde{\varsigma}$ given by (36).

*Proof:* See Appendix D. ∎

*Remark 2:* Theorem 1 provides a sufficient condition to guarantee that the quantizers will never saturate. Note that all the terms on the right sides of (16) are finite constants, which implies that the quantizers will never saturate as long as $K_x$ and $K_y$ are positive integers larger than the lower bounds in (16). In addition, (16) depends on the initial states of nodes, which is common in literature

[13], [14]. When executing the proposed algorithm in practice, we can choose $\alpha$ and $\beta$ from $(0, 1)$ arbitrarily, let $\xi$ be in close proximity to 1, and select a large enough constant and a small enough constant as the quantization level and the step size, respectively.

### B. Linear Convergence under Finite Quantization Levels

Building upon the conditions on the quantization levels in Theorem 1, the following theorem shows that the Q-DGT can linearly converge to the optimal solution at the rate of $\mathcal{O}(\xi^k)$ with $\xi \in (\hat{\rho}, 1)$.

*Theorem 2:* Suppose the conditions in Theorem 1 are satisfied. Let $\{x_i(k)\}$, $i \in \mathcal{V}$, be the sequence generated by Algorithm 1. If the quantization levels satisfy (16), then Q-DGT can linearly converge to $x^\star$ at the rate of $\mathcal{O}(\xi^k)$, i.e., $\|x_i(k) - x^\star\|_2 = \mathcal{O}(\xi^k)$ for all $i \in \mathcal{V}$.

*Proof:* Recalling (13) and (36), we can straightly obtain

$$\|\Theta(k)\|_2 \leq \|G^k\|_2\|\Theta(0)\|_2 + \tilde{\varsigma}\sum_{l=0}^{k-1}\|G^{k-1-l}\|_2\xi^l$$
$$\leq \tau\hat{\rho}^k\|\Theta(0)\|_2 + \tilde{\varsigma}\tau\sum_{l=0}^{k-1}\hat{\rho}^{k-1-l}\xi^l$$

Note that $\sum_{l=0}^{k-1}\hat{\rho}^{k-1-l}\xi^l = \xi^{k-1}\sum_{l=0}^{k-1}\left(\frac{\hat{\rho}}{\xi}\right)^{k-1-l} \leq \frac{\xi^k}{\xi-\hat{\rho}}$. Hence,

$$\|\Theta(k)\|_2 \leq \tau d_0\hat{\rho}^k + \frac{\tilde{\varsigma}\tau}{\xi-\hat{\rho}}\xi^k,$$

which implies that $\|\Theta(k)\|_2$ converges to 0 at the rate of $\mathcal{O}(\xi^k)$. Therefore, $\|\bar{x}(k) - x^\star\|_2$, $\|x(k) - \mathbf{1}_n\bar{x}(k)\|_{\mathcal{A}}$ and $\|z(k) - \pi_{\mathcal{B}}\bar{z}(k)\|_{\mathcal{B}}$ all linearly converge to 0 at the same rate. Note that

$$\|x_i(k) - x^\star\|_2 \leq \|x_i(k) - \bar{x}(k)\|_2 + \|\bar{x}(k) - x^\star\|_2$$
$$\leq \|x(k) - \mathbf{1}_n\bar{x}(k)\|_2 + \|\bar{x}(k) - x^\star\|_2,$$

which implies that $\|x_i(k) - x^\star\|_2 = \mathcal{O}(\xi^k)$ for all $i \in \mathcal{V}$ by recalling the fact that $\|x(k) - \mathbf{1}_n\bar{x}(k)\|_2 \leq \|x(k) - \mathbf{1}_n\bar{x}(k)\|_{\mathcal{A}}$. ∎

### C. 3-Level Quantization is Enough for Linear Convergence

As shown in Theorem 1, the lower bounds in (16) are finite. This inspires us to consider whether there exists a minimum number of quantization level that can preserve the linear convergence? The following theorem gives a positive answer and reveals that we can set $K_x = K_y = 1$ by appropriately tuning the associated parameters. In such an extreme scenario, each node $i \in \mathcal{V}$ can solve problem (1) with 3-level quantization.

*Theorem 3:* Suppose the conditions in Theorem 1 are satisfied. If $\alpha$ and $\beta$ are sufficiently small, then there exists $C > 0$ and $\xi \in (\hat{\rho}, 1)$ such that $K_x = K_y = 1$ is sufficient to guarantee the linear convergence of Q-DGT.

*Proof:* To prove the result, our strategy is minimizing the lower bounds obtained in (16). Particularly, if we can choose the associated parameters appropriately such that all the lower bounds in (16) can be upper bounded by 1, then it can be concluded that the quantizers will never saturate even when $K_x = K_y = 1$. In this case, Theorem 3 can be proved by recalling Theorem 2.

Now, we consider the last term of each inequality in (16). Recalling $\bar{\Upsilon} \triangleq 1 + \frac{\tilde{\varsigma}\hat{\rho}}{\xi(\xi-\hat{\rho})\|\Theta(0)\|_2} + \frac{\tilde{\varsigma}}{\xi\tau\|\Theta(0)\|_2}$ in Theorem 1 and the expression of $\tilde{\varsigma}$ in (36). If $\alpha$ and $\beta$ both tend to 0, then $\tilde{\varsigma}$ tends to 0. Since $\|\Theta(0)\|_2$, $\hat{\rho}$ and $\tau$ are all some positive constants, and $\xi$ is a constant chosen in the interval $(\hat{\rho}, 1)$, we can obtain that $\bar{\Upsilon}$ tends to 1. Therefore, the last term of each inequality in (16) can be upper bounded by $\frac{\sqrt{3}\varphi_1\tau\|\Theta(0)\|_2}{C\xi} + \frac{1}{2\xi}$ and $\frac{\sqrt{3}\varphi_2\|\Theta(0)\|_2}{C} + \frac{1}{2\xi}$, respectively. Note that $\varphi_1$ and $\varphi_2$ are constants given in (33). If we choose the constant $\xi \in (\max\{0.5, \hat{\rho}\}, 1)$, and set $C > \max\left\{\frac{2\sqrt{3}\varphi_1\tau\|\Theta(0)\|_2}{2\xi-1}, \frac{2\sqrt{3}\varphi_2\xi\tau\|\Theta(0)\|_2}{2\xi-1}\right\}$, then the last term of each
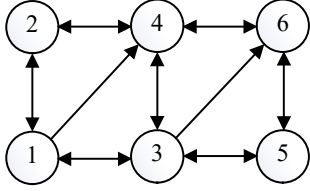
Fig. 2: The directed graph.

inequality in (16) both can be upper bounded by 1. In addition, the other terms on the right side of (16) can be upper bounded by 1 directly if we set $C > \max\{\frac{2}{3}v_1, \frac{2}{3}v_2, \frac{2\sqrt{3}\varphi_1\|\Theta(0)\|_2}{2\xi-1}\}$. In summary, there exists constants $\xi \in (\max\{0.5, \hat{\rho}\}, 1)$, and

$$C > \max\left\{\frac{2}{3}v_1, \frac{2}{3}v_2, \frac{2\sqrt{3}\varphi_1\|\Theta(0)\|_2}{2\xi-1}, \frac{2\sqrt{3}\varphi_1\tau\|\Theta(0)\|_2}{2\xi-1}, \right.$$
$$\left. \frac{2\sqrt{3}\varphi_2\xi\tau\|\Theta(0)\|_2}{2\xi-1}\right\} \quad (17)$$

such that $K_x = K_y = 1$ is sufficient to guarantee the linear convergence of Q-DGT. ∎

*Remark 3:* Quantized distributed algorithms with 3-level quantization have also been studied in [11], [12], [14], [15] to improve the communication efficiency. In contrast to the distributed optimization algorithms in [12], [14], [15], the proposed Q-DGT can achieve linear convergence. Though the quantizer has only 3 quantization levels, each node $i$ can still estimate the values of its in-neighbors $j \in \mathcal{N}_i^{in}$ iteratively via the decoding scheme (6). This can be observed from the facts that $\hat{x}_j(k) - x_j(k) = h(k)e_{x_j}(k)$ and $\hat{y}_j(k) - y_j(k) = h(k)e_{y_j}(k)$. If the quantizers never saturate, then the diminishing $h(k)$ guarantees that $\hat{x}_j(k)$ and $\hat{y}_j(k)$ tend to $x_j(k)$ and $y_j(k)$, respectively, as $k$ tends to infinity. That is why our algorithm can converge to the true solution even with 3-level quantization.

## V. NUMERICAL EXAMPLES

In this section, we apply our algorithm to the sensor fusion problem in directed networks, which has been widely adopted in the literature [25], [31]. In this problem, all sensors collectively solve the following optimization problem over the digraph decipted in Fig. 2:

$$\underset{x\in\mathbb{R}^m}{\text{minimize}} \ f(x) = \sum_{i=1}^{n} \left( \|\mathcal{M}_i x - \zeta_i\|^2 + \frac{\lambda}{2n}\|x\|^2 \right),$$

where $\mathcal{M}_i \in \mathbb{R}^{s\times m}$ and $\zeta_i \in \mathbb{R}^s$ denote the measurement matrix and the noise observation of sensor $i$, respectively, $\lambda > 0$ is the regularization parameter.

In our simulations, $\mathcal{M}_i \in \mathbb{R}^{2\times 2}$ and $\zeta_i \in \mathbb{R}^2$ are generated randomly for each $i \in \mathcal{V}$. We set $\lambda = 0.05$. $\mathcal{A}$ and $\mathcal{B}$ are designed according to the rules in Remark 2 of [25]. We first compare the convergence performance of Q-DGT with push-pull algorithm [25] under different stepsizes. The simulation results are depicted in Fig. 3(a). We can find that Q-DGT converges slower than the push-pull algorithm, which is reasonable as the performance inevitably affected by the loss of information. Despite this, the Q-DGT still maintains linear convergence, which is consistent with our theoretical results. We further compare the total cost of communicated bits between the two algorithms with $\eta = 0.008$. As shown in Fig. 3(b), the proposed Q-DGT requires less communicated bits for achieving the equal accuracy. Then, we make comparisons with the subgradient-based quantized distributed algorithms in [14] and [18]. For fair comparison, we neglect the directionality in Fig. 2 and adopt (16) for Q-DGT. The results are depicted in Fig. 4(a). It can be seen that

the convergence rate of Q-DGT outperforms that of the quantized algorithms in [14] and [18]. Finally, we verify the effectiveness of Q-DGT under different fixed numbers of quantization levels. The related parameters are chosen heuristically to meet the requirements in Theorem 3. As we can see in Fig. 4(b), the Q-DGT can still achieve linear convergence, even when the exchanged variables are respectively quantized with 3 quantization levels. In addition, a larger quantization level leads to faster convergence. This result is also reasonable since a larger quantization level implies a smaller quantization error.
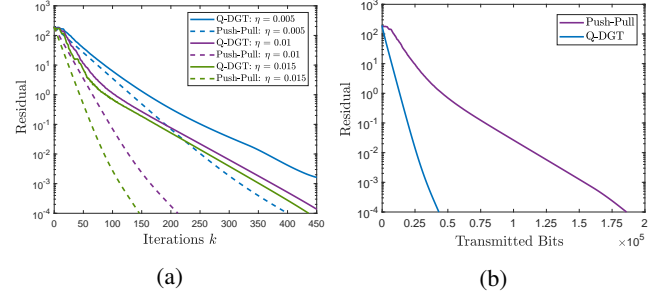


Fig. 3: (a) Comparison with push-pull in [25] under different step sizes; (b) The total communication cost of Q-DGT and push-pull.
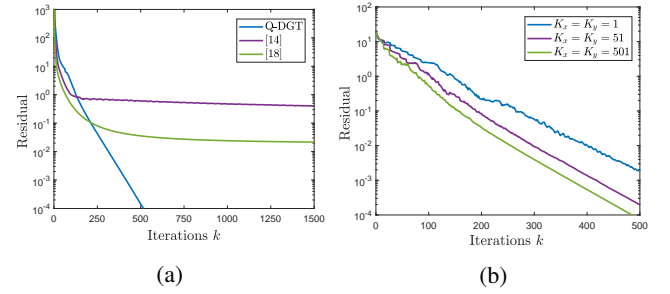


Fig. 4: (a) Comparisons with the quantized algorithms in [14] and [18]; (b) Performances of Q-DGT under different quantization levels.

## VI. CONCLUSION

In this paper, we have studied the distributed optimization problem over directed networks with quantized communications. To cope with this problem, a novel quantized distributed algorithm Q-DGT has been proposed. The lower bounds for the number of quantization levels have been explicitly derived. We have rigorously shown that Q-DGT is robust to quantization errors, and achieves linear convergence even when the exchanged variables are respectively quantized with 3 quantization levels. Future works can focus on extending the proposed algorithm to time-varying directed networks. It is also of interest to relax the conditions that preserves the convergence performance.

## APPENDIX

### A. Proof of Lemma 2.

For the clarity of presentation, we define $g(k) \triangleq \mathbf{1}_n^T \nabla F(x(k))$ and $\bar{g}(k) \triangleq \mathbf{1}_n^T \nabla F(\mathbf{1}_n\bar{x}(k))$. To prove this lemma, we first provide the following intermediate result.

*Lemma 4:* Suppose Assumptions 2-3 hold. We have $\|g(k) - \bar{g}(k)\|_2 \leq \sqrt{n}L\|x(k) - \mathbf{1}_n\bar{x}(k)\|_2$, $\|\bar{z}(k) - g(k)\|_2 \leq \|\mathbf{1}_n^T\epsilon_y(k-1)\|_2$, and $\|\bar{g}(k)\|_2 \leq nL\|\bar{x}(k) - x^\star\|_2$. If $\eta \leq \frac{1}{(\mu+L)\pi_\mathcal{A}^T\pi_\mathcal{B}}$, then

$$\|\bar{x}(k) - \eta\pi_\mathcal{A}^T\pi_\mathcal{B}\bar{g}(k) - x^\star\|_2 \leq (1 - \eta\pi_\mathcal{A}^T\pi_\mathcal{B}\mu)\|\bar{x}(k) - x^\star\|_2.$$

The first and third inequalities in Lemma 4 follow from Assumption 3 and the fact that $\|\bar{g}(k)\|_2 = \|\mathbf{1}_n^{\mathrm{T}}\nabla F(\mathbf{1}_n\bar{x}(k)) - \mathbf{1}_n^{\mathrm{T}}\nabla F(\mathbf{1}_n x^\star)\|_2$, while the second inequality can be obtained directly by applying (11). The last statement can be verified by following the similar line of Lemma 10 in [34]. Now, we begin to prove Lemma 2 by establishing the upper bounds of $\|\bar{x}(k+1) - x^\star\|_2$, $\|x(k+1) - \mathbf{1}_n\bar{x}(k+1)\|_{\mathcal{A}}$ and $\|z(k+1) - v\bar{z}(k+1)\|_{\mathcal{B}}$, respectively.

(i) In view of (10a), we have

$$
\begin{aligned}
\bar{x}(k+1) &= \bar{x}(k) - \eta\pi_{\mathcal{A}}^{\mathrm{T}}\pi_{\mathcal{B}}\bar{z}(k) - \eta\pi_{\mathcal{A}}^{\mathrm{T}}(z(k) - \pi_{\mathcal{B}}\bar{z}(k)) \\
&= \bar{x}(k) - \eta\pi_{\mathcal{A}}^{\mathrm{T}}\pi_{\mathcal{B}}\bar{g}(k) - \eta\pi_{\mathcal{A}}^{\mathrm{T}}\pi_{\mathcal{B}}(\bar{z}(k) - g(k)) \\
&\quad - \eta\pi_{\mathcal{A}}^{\mathrm{T}}\pi_{\mathcal{B}}(g(k) - \bar{g}(k)) - \eta\pi_{\mathcal{A}}^{\mathrm{T}}(z(k) - \pi_{\mathcal{B}}\bar{z}(k)) \\
&= \bar{x}(k) - \eta\pi_{\mathcal{A}}^{\mathrm{T}}\pi_{\mathcal{B}}\bar{g}(k) - \eta\pi_{\mathcal{A}}^{\mathrm{T}}\pi_{\mathcal{B}}(g(k) - \bar{g}(k)) \\
&\quad - \eta\pi_{\mathcal{A}}^{\mathrm{T}}(z(k) - \pi_{\mathcal{B}}\bar{z}(k)) - \eta\pi_{\mathcal{A}}^{\mathrm{T}}\pi_{\mathcal{B}}\mathbf{1}_n^{\mathrm{T}}\epsilon_y(k-1).
\end{aligned}
$$

Therefore, by invoking Lemma 1 and Lemma 4, we further obtain

$$
\begin{aligned}
&\|\bar{x}(k+1) - x^\star\|_2 \\
&\leq (1 - \eta\pi_{\mathcal{A}}^{\mathrm{T}}\pi_{\mathcal{B}}\mu)\|\bar{x}(k) - x^\star\|_2 + \eta\pi_{\mathcal{A}}^{\mathrm{T}}\pi_{\mathcal{B}}\|g(k) - \bar{g}(k))\|_2 \\
&\quad + \eta\|\pi_{\mathcal{A}}^{\mathrm{T}}(z(k) - \pi_{\mathcal{B}}\bar{z}(k))\|_2 + \eta\pi_{\mathcal{A}}^{\mathrm{T}}\pi_{\mathcal{B}}\left\|\mathbf{1}_n^{\mathrm{T}}\epsilon_y(k-1)\right\|_2 \\
&\leq (1 - \eta\pi_{\mathcal{A}}^{\mathrm{T}}\pi_{\mathcal{B}}\mu)\|\bar{x}(k) - x^\star\|_2 + \sqrt{n}\eta\pi_{\mathcal{A}}^{\mathrm{T}}\pi_{\mathcal{B}}\|x(k) - \mathbf{1}_n\bar{x}(k)\|_{\mathcal{A}} \\
&\quad + \eta\|z(k) - \pi_{\mathcal{B}}\bar{z}(k)\|_{\mathcal{B}} + \eta\pi_{\mathcal{A}}^{\mathrm{T}}\pi_{\mathcal{B}}\|\mathbf{1}_n^{\mathrm{T}}\epsilon_y(k-1)\|_2. \quad (18)
\end{aligned}
$$

where the fact that $\|\pi_{\mathcal{A}}\|_2 \leq 1$ has been used to obtain the last inequality.

(ii) From (9a) and (10a), along with $\mathcal{A}_\alpha\mathbf{1}_n = \mathbf{1}_n$, we obtain

$$
\begin{aligned}
&x(k+1) - \mathbf{1}_n\bar{x}(k+1) \\
&= \mathcal{A}_\alpha(x(k) - \mathbf{1}_n\bar{x}(k)) - \eta(I_n - \mathbf{1}_n\pi_{\mathcal{A}}^{\mathrm{T}})z(k) \\
&\quad + \alpha(\mathcal{A} - I_n)\sigma_x(k) \\
&= (\mathcal{A}_\alpha - \mathbf{1}_n\pi_{\mathcal{A}}^{\mathrm{T}})(x(k) - \mathbf{1}_n\bar{x}(k)) - \eta(I_n - \mathbf{1}_n\pi_{\mathcal{A}}^{\mathrm{T}})z(k) \\
&\quad + \alpha(\mathcal{A} - I_n)\sigma_x(k) \\
&= (\mathcal{A}_\alpha - \mathbf{1}_n\pi_{\mathcal{A}}^{\mathrm{T}})(x(k) - \mathbf{1}_n\bar{x}(k)) - \eta(I_n - \mathbf{1}_n\pi_{\mathcal{A}}^{\mathrm{T}})\pi_{\mathcal{B}}\bar{z}(k) \\
&\quad - \eta(I_n - \mathbf{1}_n\pi_{\mathcal{A}}^{\mathrm{T}})(z(k) - \pi_{\mathcal{B}}\bar{z}(k)) + \alpha(\mathcal{A} - I_n)\sigma_x(k),
\end{aligned}
$$

where the fact that $\mathbf{1}_n\pi_{\mathcal{A}}^{\mathrm{T}}(x(k) - \mathbf{1}_n\bar{x}(k)) = \mathbf{0}_n$ has been exploited to obtain the second equality. By employing Lemma 1, we obtain

$$
\begin{aligned}
&\|x(k+1) - \mathbf{1}_n\bar{x}(k+1)\|_{\mathcal{A}} \\
&\leq \sigma_{\mathcal{A}}\|x(k) - \mathbf{1}_n\bar{x}(k)\|_{\mathcal{A}} + \alpha\|(\mathcal{A} - I_n)\sigma_x(k)\|_{\mathcal{A}} \\
&\quad + \eta\kappa_1\|z(k) - \pi_{\mathcal{B}}\bar{z}(k)\|_{\mathcal{A}} + \eta\kappa_1\kappa_2\delta_{\mathcal{A}2}\|\bar{z}(k)\|_2. \quad (19)
\end{aligned}
$$

Now, it remains to establish an upper bound for $\|\bar{z}(k)\|_2$. Note that

$$
\begin{aligned}
\|\bar{z}(k)\|_2 &\leq \|\bar{z}(k) - g(k)\|_2 + \|g(k) - \bar{g}(k)\|_2 + \|\bar{g}(k)\|_2 \\
&\leq \|\mathbf{1}_n^{\mathrm{T}}\epsilon_y(k-1)\|_2 + \sqrt{n}L\|x(k) - \mathbf{1}_n\bar{x}(k)\|_{\mathcal{A}} \\
&\quad + nL\|\bar{x}(k) - x^\star\|_2. \quad (20)
\end{aligned}
$$

By substituting (20) into (19), we obtain

$$
\begin{aligned}
&\|x(k+1) - \mathbf{1}_n\bar{x}(k+1)\|_{\mathcal{A}} \\
&\leq \left(\sigma_{\mathcal{A}} + \sqrt{n}\eta L\kappa_1\kappa_2\delta_{\mathcal{A}2}\right)\|x(k) - \mathbf{1}_n\bar{x}(k)\|_{\mathcal{A}} \\
&\quad + \eta\delta_{\mathcal{A}\mathcal{B}}\kappa_1\|z(k) - \pi_{\mathcal{B}}\bar{z}(k)\|_{\mathcal{B}} + n\eta L\kappa_1\kappa_2\delta_{\mathcal{A}2}\|\bar{x}(k) - x^\star\|_2 \\
&\quad + \eta\kappa_1\kappa_2\delta_{\mathcal{A}2}\|\mathbf{1}_n^{\mathrm{T}}\epsilon_y(k-1)\|_2 + \alpha\|(\mathcal{A} - I_n)\sigma_x(k)\|_{\mathcal{A}}. \quad (21)
\end{aligned}
$$

(iii) In light of relations (9b) and (10b), we have

$$
\begin{aligned}
&z(k+1) - \pi_{\mathcal{B}}\bar{z}(k+1) \\
&= (\mathcal{B}_\beta - \pi_{\mathcal{B}}\mathbf{1}_n^{\mathrm{T}})(z(k) - \pi_{\mathcal{B}}\bar{z}(k)) \\
&\quad + (I_n - \pi_{\mathcal{B}}\mathbf{1}_n^{\mathrm{T}})(\nabla F(x(k+1)) - \nabla F(x(k))) \\
&\quad + (I_n - \pi_{\mathcal{B}}\mathbf{1}_n^{\mathrm{T}})(\epsilon_y(k) - \epsilon_y(k-1)),
\end{aligned}
$$

where the equality follows from the definition of $\mathcal{B}_\beta$ and the fact that $\pi_{\mathcal{B}}\mathbf{1}_n^{\mathrm{T}}\pi_{\mathcal{B}} = \pi_{\mathcal{B}}$. Hence, we obtain

$$
\begin{aligned}
&\|z(k+1) - \pi_{\mathcal{B}}\bar{z}(k+1)\|_{\mathcal{B}} \\
&\leq \delta_{\mathcal{B}2}\|I_n - \pi_{\mathcal{B}}\mathbf{1}_n^{\mathrm{T}}\|_{\mathcal{B}}\|\nabla F(x(k+1)) - \nabla F(x(k))\|_2 \\
&\quad + \|I_n - \pi_{\mathcal{B}}\mathbf{1}_n^{\mathrm{T}}\|_{\mathcal{B}}\|\epsilon_y(k) - \epsilon_y(k-1)\|_{\mathcal{B}} \\
&\quad + \sigma_{\mathcal{B}}\|z(k) - \pi_{\mathcal{B}}\bar{z}(k)\|_{\mathcal{B}}, \quad (22)
\end{aligned}
$$

where Lemma 1 has been utilized to obtain the above inequality. Now, it remains to bound $\|\nabla F(x(k+1)) - \nabla F(x(k))\|_2$. Note that

$$
\begin{aligned}
&\|\nabla F(x(k+1)) - \nabla F(x(k))\|_2 \\
&\leq L\|x(k+1) - x(k)\|_2 \\
&= L\|\mathcal{A}_\alpha x(k) - x(k) + \alpha(\mathcal{A} - I_n)\sigma_x(k) - \eta z(k)\|_2 \\
&\leq L\|\mathcal{A}_\alpha - I_n\|_2\|x(k) - \mathbf{1}_n\bar{x}(k)\|_2 + \eta L\|z(k) - \pi_{\mathcal{B}}\bar{z}(k)\|_2 \\
&\quad + \eta L\|\bar{z}(k)\|_2 + \alpha L\|(\mathcal{A} - I_n)\sigma_x(k)\|_2. \quad (23)
\end{aligned}
$$

where the fact that $\|\pi_{\mathcal{B}}\|_2 \leq 1$ has been used to obtain the last inequality. Then, by substituting (20) and (23) into (22), we can obtain

$$
\begin{aligned}
&\|z(k+1) - \pi_{\mathcal{B}}\bar{z}(k+1)\|_{\mathcal{B}} \\
&\leq \delta_{\mathcal{B}2}\kappa_3\left(L\kappa_4 + \sqrt{n}\eta L^2\right)\|x(k) - \mathbf{1}_n\bar{x}(k)\|_{\mathcal{A}} \\
&\quad + \delta_{\mathcal{B}2}\alpha L\kappa_3\|(\mathcal{A} - I_n)\sigma_x(k)\|_2 + \delta_{\mathcal{B}2}\eta L\kappa_3\|\mathbf{1}_n^{\mathrm{T}}\epsilon_y(k-1)\|_2 \\
&\quad + \kappa_3\|\epsilon_y(k) - \epsilon_y(k-1)\|_{\mathcal{B}} + \delta_{\mathcal{B}2}n\eta L^2\kappa_3\|\bar{x}(k) - x^\star\|_2 \\
&\quad + (\sigma_{\mathcal{B}} + \delta_{\mathcal{B}2}\eta L\kappa_3)\|z(k) - \pi_{\mathcal{B}}\bar{z}(k)\|_{\mathcal{B}}. \quad (24)
\end{aligned}
$$

Combining (18), (21) and (24), we can obtain (13) with

$$
G = \begin{bmatrix}
1 - \eta\pi_{\mathcal{A}}^{\mathrm{T}}\pi_{\mathcal{B}}\mu & \sqrt{n}\eta\pi_{\mathcal{A}}^{\mathrm{T}}\pi_{\mathcal{B}} & \eta \\
n\eta L\kappa_1\kappa_2\delta_{\mathcal{A}2} & \sigma_{\mathcal{A}} + \sqrt{n}\eta L\kappa_1\kappa_2\delta_{\mathcal{A}2} & \eta\kappa_1\delta_{\mathcal{A}\mathcal{B}} \\
n\eta L^2\kappa_3\delta_{\mathcal{B}2} & \kappa_3(L\kappa_4 + \sqrt{n}\eta L^2)\delta_{\mathcal{B}2} & \sigma_{\mathcal{B}} + \eta L\kappa_3\delta_{\mathcal{B}2}
\end{bmatrix}
$$

$$(25)$$

and $\varsigma(k) = (\varsigma_1(k), \varsigma_2(k), \varsigma_3(k))^{\mathrm{T}}$ given by

$$
\begin{aligned}
\varsigma_1(k) &= \eta\pi_{\mathcal{A}}^{\mathrm{T}}\pi_{\mathcal{B}}\|\mathbf{1}_n^{\mathrm{T}}\epsilon_y(k-1)\|_2, \\
\varsigma_2(k) &= \eta\kappa_1\kappa_2\delta_{\mathcal{A}2}\|\mathbf{1}_n^{\mathrm{T}}\epsilon_y(k-1)\|_2 + \alpha\|(\mathcal{A} - I_n)\sigma_x(k)\|_{\mathcal{A}}, \\
\varsigma_3(k) &= \kappa_3\|\epsilon_y(k) - \epsilon_y(k-1)\|_{\mathcal{B}} + \eta L\kappa_3\delta_{\mathcal{B}2}\|\mathbf{1}_n^{\mathrm{T}}\epsilon_y(k-1)\|_2 \\
&\quad + \alpha L\kappa_3\delta_{\mathcal{B}2}\|(\mathcal{A} - I_n)\sigma_x(k)\|_2, \quad (26)
\end{aligned}
$$

which completes the proof. ∎

### B. Proof of Lemma 3.

To achieve this goal, we need to provide a sufficient condition under which $G_{ii} < 1$ and $\det(I - G) > 0$ can be guaranteed [25]. We first ensure that $G_{ii} < 1$ hold for $i = 1, 2, 3$. Clearly, if we set $\eta \leq \frac{1}{(\mu+L)\pi_{\mathcal{A}}^{\mathrm{T}}\pi_{\mathcal{B}}}$, then $0 < G_{11} < 1$. We can also verify that if $\eta \leq \min\{\frac{1-\sigma_{\mathcal{A}}}{2\sqrt{n}L\kappa_1\kappa_2\delta_{\mathcal{A}2}}, \frac{1-\sigma_{\mathcal{B}}}{2L\kappa_3\delta_{\mathcal{B}2}}\}$, then $G_{22} < 1$ and $G_{33} < 1$ both hold. Now we turn our attention to $\det(I_3 - G)$. Note that

$$
\begin{aligned}
&\det(I_3 - G) \\
&= \eta\pi_{\mathcal{A}}^{\mathrm{T}}\pi_{\mathcal{B}}\mu[1 - (\sigma_{\mathcal{A}} + \sqrt{n}\eta L\kappa_1\kappa_2\delta_{\mathcal{A}2})][1 - (\sigma_{\mathcal{B}} + \eta L\kappa_3\delta_{\mathcal{B}2})] \\
&\quad - \eta\pi_{\mathcal{A}}^{\mathrm{T}}\pi_{\mathcal{B}}\mu(\eta\kappa_1\delta_{\mathcal{A}\mathcal{B}})[\kappa_3(L\kappa_4 + \sqrt{n}\eta L^2)\delta_{\mathcal{B}2}] \\
&\quad - \sqrt{n}\eta\eta\pi_{\mathcal{A}}^{\mathrm{T}}\pi_{\mathcal{B}}(n\eta L\kappa_1\kappa_2\delta_{\mathcal{A}2})[1 - (\sigma_{\mathcal{B}} + \eta L\kappa_3\delta_{\mathcal{B}2})] \\
&\quad - \sqrt{n}\eta\eta\pi_{\mathcal{A}}^{\mathrm{T}}\pi_{\mathcal{B}}(\eta\kappa_1\delta_{\mathcal{A}\mathcal{B}})(n\eta L^2\kappa_3\delta_{\mathcal{B}2}) \\
&\quad - \eta(n\eta L\kappa_1\kappa_2\delta_{\mathcal{A}2})[\kappa_3(L\kappa_4 + \sqrt{n}\eta L^2)\delta_{\mathcal{B}2}] \\
&\quad - \eta[1 - (\sigma_{\mathcal{A}} + \sqrt{n}\eta L\kappa_1\kappa_2\delta_{\mathcal{A}2})](n\eta L^2\kappa_3\delta_{\mathcal{B}2}).
\end{aligned}
$$

In light of $\eta \leq \min\{\frac{1-\sigma_{\mathcal{A}}}{2\sqrt{n}L\kappa_1\kappa_2\delta_{\mathcal{A}2}}, \frac{1-\sigma_{\mathcal{B}}}{2L\kappa_3\delta_{\mathcal{B}2}}\}$, we have $\frac{1-\sigma_{\mathcal{A}}}{2} \leq 1 - (\sigma_{\mathcal{A}} + \sqrt{n}\eta L\kappa_1\kappa_2\delta_{\mathcal{A}2}) \leq 1 - \sigma_{\mathcal{A}}$ and $\frac{1-\sigma_{\mathcal{B}}}{2} \leq 1 - (\sigma_{\mathcal{B}} + \delta_{\mathcal{B}2}\eta L\kappa_3) \leq 1 - \sigma_{\mathcal{B}}$. Therefore, a sufficient condition for $\det(I - G) > 0$ is

$$
\frac{1}{4}\eta\pi_{\mathcal{A}}^{\mathrm{T}}\pi_{\mathcal{B}}\mu(1 - \sigma_{\mathcal{A}})(1 - \sigma_{\mathcal{B}}) - \eta(1 - \sigma_{\mathcal{A}})(n\eta L^2\kappa_3\delta_{\mathcal{B}2})
$$

$$-\eta\pi_{\mathcal{A}}^{\mathrm{T}}\pi_{\mathcal{B}}\mu(\eta\kappa_1\delta_{\mathcal{AB}})[\kappa_3(L\kappa_4+\sqrt{n}\eta L^2)\delta_{\mathcal{B}2}]$$
$$-\sqrt{n}\eta\eta\pi_{\mathcal{A}}^{\mathrm{T}}\pi_{\mathcal{B}}(n\eta L\kappa_1\kappa_2\delta_{\mathcal{A}2})(1-\sigma_{\mathcal{B}})$$
$$-\sqrt{n}\eta\eta\pi_{\mathcal{A}}^{\mathrm{T}}\pi_{\mathcal{B}}(\eta\kappa_1\delta_{\mathcal{AB}})(n\eta L^2\kappa_3\delta_{\mathcal{B}2})$$
$$-\eta(n\eta L\kappa_1\kappa_2\delta_{\mathcal{A}2})[\kappa_3(L\kappa_4+\sqrt{n}\eta L^2)\delta_{\mathcal{B}2}]>0. \quad (27)$$

Now, the inequality (27) can be rewritten as $\Gamma_1\eta^2+\Gamma_2\eta-\Gamma_3<0$ with $\Gamma_i$, $i=1,2,3$, given by

$$\Gamma_1 = \sqrt{n}\kappa_1\kappa_3L^2\delta_{\mathcal{B}2}[(n+\mu)\pi_{\mathcal{A}}^{\mathrm{T}}\pi_{\mathcal{B}}\delta_{\mathcal{AB}}+n\kappa_2L\delta_{\mathcal{A}2}],$$
$$\Gamma_2 = \kappa_1L\pi_{\mathcal{A}}^{\mathrm{T}}\pi_{\mathcal{B}}[n^{\frac{3}{2}}\kappa_2\delta_{\mathcal{A}2}(1-\sigma_{\mathcal{B}})+\mu\kappa_3\kappa_4\delta_{\mathcal{AB}}\delta_{\mathcal{B}2}]$$
$$+n\kappa_3L^2\delta_{\mathcal{B}2}[1-\sigma_{\mathcal{A}}+\kappa_1\kappa_2\kappa_4\delta_{\mathcal{A}2}],$$
$$\Gamma_3 = \frac{1}{4}\mu\pi_{\mathcal{A}}^{\mathrm{T}}\pi_{\mathcal{B}}(1-\sigma_{\mathcal{A}})(1-\sigma_{\mathcal{B}}), \quad (28)$$

Therefore, it can be derived that $\eta\leq\frac{2\Gamma_3}{\Gamma_2+\sqrt{\Gamma_2^2+4\Gamma_1\Gamma_3}}$, which completes the proof. ∎

### C. Proof of Theorem 1.

To ensure that the finite-level quantizers never saturate, the scaled "innovation" $\frac{1}{h(k)}(x_j(k)-\hat{x}_j(k-1))$ and $\frac{1}{h(k)}(y_j(k)-\hat{y}_j(k-1))$ must lie in a bounded region. To achieve the goal, we first establish the upper bounds for $\frac{1}{h(k)}\|x_i(k)-\hat{x}_i(k-1)\|_\infty$ and $\frac{1}{h(k)}\|y_i(k)-\hat{y}_i(k-1)\|_\infty$, respectively. Then, the obtained upper bounds lead us to propose an update rule of the quantization levels, under which we prove the unsaturation of quantizers by mathematical induction. Finaly, we show that (16) suffices for the given update rule.

**Step 1: Bound $\|x_i(k)-\hat{x}_i(k-1)\|_\infty$ and $\|y_i(k)-\hat{y}_i(k-1)\|_\infty$.**
Let $e_{x_i}(k)\triangleq Q_{K_x}(\frac{x_i(k)-\hat{x}_i(k-1)}{h(k)})-\frac{x_i(k)-\hat{x}_i(k-1)}{h(k)}$ and $e_{y_i}(k)\triangleq Q_{K_y}(\frac{y_i(k)-\hat{y}_i(k-1)}{h(k)})-\frac{y_i(k)-\hat{y}_i(k-1)}{h(k)}$. Recalling (5) and (6), we can obtain $\hat{x}_j(k)=x_j(k)+h(k)e_{x_j}(k)$. Then

$$\|x_i(k)-\hat{x}_i(k-1)\|_\infty$$
$$\leq\|x_i(k)-x_i(k-1)\|_\infty+h(k-1)\|e_{x_i}(k-1)\|_\infty. (29)$$

The first term on the right side of (29) can be further calculated as

$$\|x_i(k)-x_i(k-1)\|_\infty\leq\alpha\sum_{j=1}^n\|x_i(k-1)-x_j(k-1)\|_\infty$$
$$+\eta\|z_i(k-1)\|_\infty+\alpha\sum_{j=1}^n\|\sigma_{x_i}(k-1)-\sigma_{x_j}(k-1)\|_\infty, (30)$$

where the inequality follows from (9a), the definition of $\mathcal{A}_\alpha$, and the row stochasticity of $\mathcal{A}$. In the following, we will establish the upper bounds for the three terms on the right side of (30), respectively.

For the first term, it can be calculated as follows

$$\sum_{j=1}^n\|x_i(k-1)-x_j(k-1)\|_\infty \leq \sqrt{2}(n+\frac{1}{2})\Theta_2(k-1), (31)$$

where the Jensen's inequality and the facts that $\|x_i(k-1)-x_j(k-1)\|_\infty^2\leq 2\|x_i(k-1)-\bar{x}(k-1)\|_\infty^2+2\|\bar{x}(k-1)-x_j(k-1)\|_\infty^2$ and $\|x_i(k-1)-\bar{x}(k-1)\|_\infty^2\leq\|x_i(k-1)-\bar{x}(k-1)\|_2^2$ have been exploited to obtain the above inequality.

For the second term on the right side of (30), we have

$$\|z_i(k-1)\|_\infty$$
$$\leq\|z(k-1)-\pi_{\mathcal{B}}\bar{z}(k-1)\|_\infty+\|\pi_{\mathcal{B}}\|_\infty\|\bar{g}(k-1)\|_\infty$$
$$+\|\pi_{\mathcal{B}}\|_\infty\|\bar{z}(k-1)-g(k-1)\|_\infty$$
$$+\|\pi_{\mathcal{B}}\|_\infty\|g(k-1)-\bar{g}(k-1)\|_\infty$$
$$\leq\Theta_3(k-1)+\sqrt{n}L\Theta_2(k-1)+nL\Theta_1(k-1)$$
$$+n\beta h(k-2)\max_{i\in\mathcal{V}}\|e_{y_i}(k-2)\|_\infty, (32)$$

where Lemma 1 and Lemma 4 have been employed to obtain the above inequality.

It only remains to bound the last term in (30). By using the fact that $\hat{x}_j(k)=x_j(k)+h(k)e_{x_j}(k)$ again, we obtain $\sum_{j=1}^n\|\sigma_{x_i}(k-1)-\sigma_{x_j}(k-1)\|_\infty\leq 2nh(k-1)\max_{i\in\mathcal{V}}\|e_{x_i}(k-1)\|_\infty$.
Define

$$\varphi_1 \triangleq \max\left\{\sqrt{2}(n+\frac{1}{2})\alpha+\eta\sqrt{n}L,\eta,\eta nL\right\},$$
$$\varphi_2 \triangleq \max\left\{1,\sqrt{n}L,nL\right\}. (33)$$

Combining the above inequalities, we can obtain

$$\|x_i(k)-\hat{x}_i(k-1)\|_\infty\leq\sqrt{3}\varphi_1\|\Theta(k-1)\|_2$$
$$+(2\alpha n+1)h(k-1)\max_{i\in\mathcal{V}}\|e_{x_i}(k-1)\|_\infty$$
$$+n\eta\beta h(k-2)\max_{i\in\mathcal{V}}\|e_{y_i}(k-2)\|_\infty. (34)$$

From (34), we can observe that if the quantizers never saturate, then $x_j(k)-\hat{x}_j(k-1)$ will decay to zero at the speed of the same order of $h(k)$ since $h(k-1)=\frac{C}{\xi}\xi^k$. Following the similar line above, we can further obtain

$$\|y_i(k)-\hat{y}_i(k-1)\|_\infty\leq\sqrt{3}\varphi_2\|\Theta(k)\|_2$$
$$+(n\beta+1)h(k-1)\max_{i\in\mathcal{V}}\|e_{y_i}(k-1)\|_\infty.$$

**Step 2: Demonstrate the unsaturation.**

In this part, we first consider the following update rule of the quantization levels instead

$$K_x(0)\geq\frac{v_1}{C}-\frac{1}{2}, \ K_y(0)\geq\frac{v_2}{C}-\frac{1}{2}$$
$$K_x(1)\geq\frac{\sqrt{3}\varphi_1\|\Theta(0)\|_2}{C\xi}+\frac{2\alpha n+1}{2\xi}-\frac{1}{2}$$
$$K_x(k)\geq\frac{\sqrt{3}\varphi_1\tau\|\Theta(0)\|_2}{C\xi}\Upsilon_1(k)+\frac{2\alpha n+1}{2\xi}+\frac{n\eta\beta}{2\xi^2}-\frac{1}{2},k\geq 2$$
$$K_y(k)\geq\frac{\sqrt{3}\varphi_2\tau\|\Theta(0)\|_2}{C}\Upsilon_2(k)+\frac{n\beta+1}{2\xi}-\frac{1}{2}, \ k\geq 1 \quad (35)$$

where $\Upsilon_1(k)=(\frac{\hat{\rho}}{\xi})^{k-1}+\frac{\check{\varsigma}}{\xi\|\Theta(0)\|_2}\sum_{l=0}^{k-3}(\frac{\hat{\rho}}{\xi})^{k-2-l}+\frac{\check{\varsigma}}{\xi\tau\|\Theta(0)\|_2}$ and $\Upsilon_2(k)=(\frac{\hat{\rho}}{\xi})^k+\frac{\check{\varsigma}}{\xi\|\Theta(0)\|_2}\sum_{l=0}^{k-2}(\frac{\hat{\rho}}{\xi})^{k-1-l}+\frac{\check{\varsigma}}{\xi\tau\|\Theta(0)\|_2}$.

Now, we show the unsaturation of the quantizers under the rule (35) by mathematical induction. Considering the case $k=0$, we have $\frac{\|x_i(k)-\hat{x}_i(k-1)\|_\infty}{h(k)}\leq\frac{\|x_i(0)\|_\infty}{C}\leq K_x(0)+\frac{1}{2}$ and $\frac{\|y_i(k)-\hat{y}_i(k-1)\|_\infty}{h(k)}\leq\frac{\|y_i(0)\|_\infty}{C}\leq K_y(0)+\frac{1}{2}$, which indicates that the quantizers are not saturated for $k=0$. Therefore, $\max_{i\in\mathcal{V}}\|e_{x_i}(0)\|_\infty\leq\frac{1}{2}$ and $\max_{i\in\mathcal{V}}\|e_{y_i}(0)\|_\infty\leq\frac{1}{2}$ both hold, which further can be exploited to calculate the upper bounds of $\varsigma_i(0)$ via (26), denoted by $\bar{\varsigma}_i(0)$, for $i=1,2,3$. Define $\hat{\varsigma}(0)\triangleq\|\bar{\varsigma}(0)\|_2$ with $\bar{\varsigma}(0)=(\bar{\varsigma}_1(0),\bar{\varsigma}_2(0),\bar{\varsigma}_3(0))^{\mathrm{T}}$. Recalling (13), we can obtain $\|\Theta(1)\|_2\leq\tau\hat{\rho}\|\Theta(0)\|_2+\hat{\varsigma}(0)$.

Now, considering the case $k=1$. From (34), we can obtain $\frac{\|x_i(1)-\hat{x}_i(0)\|_\infty}{h(1)}\leq\frac{\sqrt{3}\varphi_1}{C\xi}\|\Theta(0)\|_2+\frac{2\alpha n+1}{\xi}\max_{i\in\mathcal{V}}\|e_{x_i}(0)\|_\infty\leq\frac{\sqrt{3}\varphi_1}{C\xi}\|\Theta(0)\|_2+\frac{2\alpha n+1}{2\xi}\leq K_x(1)+\frac{1}{2}$. Similarly, it can be easily verified that $\frac{\|y_i(1)-\hat{y}_i(0)\|_\infty}{h(1)}\leq K_y(1)+\frac{1}{2}$. These two inequalities imply that the quantizers are not saturated at $k=1$ as well. Then, we have $\max_{i\in\mathcal{V}}\|e_{x_i}(\nu)\|_\infty\leq\frac{1}{2}$ and $\max_{i\in\mathcal{V}}\|e_{y_i}(\nu)\|_\infty\leq\frac{1}{2}$ for $\nu\in\{0,1\}$, which further can be utilized to compute $\hat{\varsigma}(1)$. Hence, we can obtain $\|\Theta(2)\|_2\leq\tau\|\Theta(0)\|_2\hat{\rho}^2+\tau\hat{\rho}\hat{\varsigma}(0)+\hat{\varsigma}(1)$.

From the above observations, it can be seen that our basic idea is to exploit the non-saturation property at each step, i.e., $\max_{i\in\mathcal{V}}\|e_{x_i}(\nu)\|_\infty\leq\frac{1}{2}$ and $\max_{i\in\mathcal{V}}\|e_{y_j}(\nu)\|_\infty\leq\frac{1}{2}$ for $\nu\in\{0,1,...,k-1\}$, then we can derive the upper bounds of $\|\varsigma(\nu)\|_2$. In this way, the upper bounds of $\|\Theta(\nu+1)\|_2$ can be obtained, which further helps us to derive the non-saturation condition at step $k$. In

other words, if the quantizers are not saturated for all $k \leq k'$, we can obtain $\hat{\varsigma}(1), ..., \hat{\varsigma}(k')$ with

$$\hat{\varsigma}(l) = \|(\bar{\varsigma}_1(l), \bar{\varsigma}_2(l), \bar{\varsigma}_3(l))^{\mathrm{T}}\|_2, \quad l \leq k',$$

where $\bar{\varsigma}(l) = \xi^l \bar{\varsigma}$ and the elements of the vector $\bar{\varsigma} \in \mathbb{R}^3$ is given by: $\bar{\varsigma}_1 = \frac{1}{2\xi}\eta\pi_{\mathcal{A}}^{\mathrm{T}}\pi_{\mathcal{B}}n\sqrt{m}\beta C$, $\bar{\varsigma}_2 = \frac{1}{2\xi}\eta\kappa_1\kappa_2\delta_{\mathcal{A}2}n\sqrt{m}\beta C + \frac{\alpha}{2}\sqrt{mn}\delta_{\mathcal{A}2}\kappa_4 C$, $\bar{\varsigma}_3 = \frac{1}{2\xi}\delta_{\mathcal{B}2}\kappa_3 n\sqrt{m}\beta C(1 + \xi + \eta L) + \frac{1}{2}\alpha\delta_{\mathcal{B}2}\kappa_3\kappa_4 L\sqrt{mn}C$. Note that each element of the vector $\bar{\varsigma}$ is a finite constant. We further define the constant $\tilde{\varsigma}$ by:

$$\tilde{\varsigma} \triangleq \|(\bar{\varsigma}_1, \bar{\varsigma}_2, \bar{\varsigma}_3)^{\mathrm{T}}\|_2. \tag{36}$$

Then, we obtain that $\|\Theta(\iota)\|_2 \leq \|G\|_2^\iota\|\Theta(0)\|_2 + \tilde{\varsigma}\sum_{l=0}^{\iota-1}\|G\|_2^{\iota-1-l}\xi^l \leq \tau\|\Theta(0)\|_2\hat{\rho}^\iota + \tilde{\varsigma}\tau\sum_{l=0}^{\iota-2}\hat{\rho}^{\iota-1-l}\xi^l + \tilde{\varsigma}\xi^{\iota-1}$ holds, for $\iota \in \{2, 3..., k'+1\}$.

Considering the case $k = k' + 1$ ($k' \geq 2$). From (34), we have

$$\frac{\|x_i(k) - \hat{x}_i(k-1)\|_\infty}{h(k)}$$
$$\leq \frac{\sqrt{3}\varphi_1}{C\xi^k}\left(\tau\|\Theta(0)\|_2\hat{\rho}^{k-1} + \tilde{\varsigma}\tau\sum_{l=0}^{k-3}\hat{\rho}^{k-2-l}\xi^l + \tilde{\varsigma}\xi^{k-2}\right)$$
$$+ \frac{2\alpha n + 1}{2\xi} + \frac{n\eta\beta}{2\xi^2} \leq K_x(k) + \frac{1}{2}. \tag{37}$$

Similarly, with some tedious calculations, it can also be concluded that $\frac{\|y_i(k) - \hat{y}_i(k-1)\|_\infty}{h(k-1)} \leq K_y(k) + \frac{1}{2}$. In summary, the quantizers will never saturate under the rule (35). Recalling $\Upsilon_1(k)$ and $\Upsilon_2(k)$ in (35), it can be verified that they both can be upper bounded by

$$\bar{\Upsilon} = 1 + \frac{\tilde{\varsigma}\hat{\rho}}{\xi(\xi - \hat{\rho})\|\Theta(0)\|_2} + \frac{\tilde{\varsigma}}{\xi\tau\|\Theta(0)\|_2}. \tag{38}$$

Note that $\|\Theta(0)\|_2$, $\hat{\rho}$ and $\tau$ are all some positive constants, $\tilde{\varsigma}$ is a positive constant given in (36), and $\xi$ is a constant chosen in the interval $(\hat{\rho}, 1)$. Hence, $\bar{\Upsilon}$ is a constant, and (16) suffices for the update rule (35), which completes the proof. ∎

## REFERENCES

[1] K. You, and L. Xie, "Network topology and communication data rate for consensusability of discrete-time multi-agent systems," *IEEE Trans. Autom. Control*, vol. 56, no. 10, pp. 2262–2275, 2011.

[2] J. Zhang, K. You, and K. Cai, "Distributed dual gradient tracking for resource allocation in unbalanced networks," *IEEE Trans. Signal Process.*, vol. 68, pp. 2186–2198, 2020.

[3] X. Li, X. Yi, and L. Xie, "Distributed online optimization for multi-agent networks with coupled inequality constraints," *IEEE Trans. Autom. Control*, 2020.

[4] S. Yuan, H. Wang, and L. Xie, "Survey on localization systems and algorithms for unmanned systems," *Unmanned Syst.*, vol. 9, no. 2, pp. 129–163, 2021.

[5] M. Ye, G. Hu, F. L. Lewis, and L. Xie, "A unified strategy for solution seeking in graphical N-coalition noncooperative games," *IEEE Trans. Autom. Control*, vol. 64, no. 11, pp. 4645–4652, 2019.

[6] A. Nedić, A. Olshevsky, and M. G. Rabbat, "Network topology and communication-computation tradeoffs in decentralized optimization," *Proceedings of the IEEE*, vol. 106, no. 5, pp. 953–976, 2018.

[7] A. Nedić, "Distributed gradient methods for convex machine learning problems in networks: Distributed optimization," *IEEE Signal Process. Mag.*, vol. 37, no. 3, pp. 92–101, 2020.

[8] S. Magnússon, C. Enyioha, N. Li, C. Fischione, and V. Tarokh, "Convergence of limited communication gradient methods," *IEEE Trans. Autom. Control*, vol. 63, no. 5, pp. 1356–1371, 2018.

[9] T. T. Doan, S. T. Maguluri, and J. Romberg, "Convergence rates of distributed gradient methods under random quantization: A stochastic approximation approach," *IEEE Trans. Autom. Control*, 2020.

[10] M. Doostmohammadian, A. Aghasi, M. Pirani, E. Nekouei, U. A. Khan, and T. Charalambous, "Fast-convergent anytime-feasible dynamics for distributed allocation of resources over switching sparse networks with quantized communication links," in *Eur. Control Conf.*, 2022, pp. 84–89.

[11] T. Li, M. Fu, L. Xie, and J. Zhang, "Distributed consensus with limited communication data rate," *IEEE Trans. Autom. Control*, vol. 56, no. 2, pp. 279–292, 2010.

[12] J. Lei, P. Yi, G. Shi, and B. D. Anderson, "Distributed algorithms with finite data rates that solve linear equations," *SIAM J. Optim.*, vol. 30, no. 2, pp. 1191–1222, 2020.

[13] H. Li, C. Huang, Z. Wang, G. Chen, and H. G. Ahmad Umar, "Computation-efficient distributed algorithm for convex optimization over time-varying networks with limited bandwidth communication," *IEEE Trans. Signal Inf. Process. over Netw.*, vol. 6, pp. 140–151, 2020.

[14] P. Yi, and Y. Hong, "Quantized subgradient algorithm and data-rate analysis for distributed optimization," *IEEE Trans. Control Netw. Syst.*, vol. 1, no. 4, pp. 380–392, 2014.

[15] J. Zhang, K. You, and T. Başar, "Distributed discrete-time optimization in multi-agent networks using only sign of relative state," *IEEE Trans. Autom. Control*, vol. 64, no. 6, pp. 2352–2367, 2019.

[16] A. Reisizadeh, A. Mokhtari, H. Hassani, and R. Pedarsani, "An exact quantized decentralized gradient descent algorithm," *IEEE Trans. Signal Process.*, vol. 67, no. 19, pp. 4934–4947, 2019.

[17] A. Nedić, and A. Ozdaglar, "Distributed subgradient methods for multi-agent optimization," *IEEE Trans. Autom. Control*, vol. 54, no. 1, pp. 48–61, 2009.

[18] H. Taheri, A. Mokhtari, H. Hassani, and R. Pedarsani, "Quantized decentralized stochastic learning over directed graphs," in *Int. Conf. Mach. Learn. (ICML)*, 2020, pp. 9324–9333.

[19] D. Kovalev, A. Koloskova, M. Jaggi, P. Richtarik, and S. U. Stich, "A linearly convergent algorithm for decentralized optimization: Sending less bits for free!" in *Int. Conf. Artif. Intell. Statist. (AISTATS)*, 2021, pp. 4087–4095.

[20] X. Liu, Y. Li, R. Wang, J. Tang, and M. Yan, "Linear convergent decentralized optimization with compression," in *Int. Conf. Learn. Repres. (ICLR)*, 2021.

[21] C. Lee, N. Michelusi, and G. Scutari, "Finite rate quantized distributed optimization with geometric convergence," in *Proc. 52nd Asilomar Conf. Signals, Syst., Comput.*, 2018, pp. 1876–1880.

[22] S. Magnússon, H. Shokri-Ghadikolaei, and N. Li. "On maintaining linear convergence of distributed learning and optimization under limited communication," *IEEE Trans. Signal Process.*, vol. 68, pp. 6101–6116, 2020.

[23] P. Xie, K. You, R. Tempo, S. Song, and C. Wu, "Distributed convex optimization with inequality constraints over time-varying unbalanced digraphs," *IEEE Trans. Autom. Control*, vol. 63, no. 12, pp. 4331–4337, 2018.

[24] A. Nedić, and A. Olshevsky, "Distributed optimization over time-varying directed graphs," *IEEE Trans. Autom. Control*, vol. 60, no. 3, pp. 601–615, 2015.

[25] S. Pu, W. Shi, J. Xu, and A. Nedić, "Push-pull gradient methods for distributed optimization in networks," *IEEE Trans. Autom. Control*, 2020.

[26] R. Xin, and U. A. Khan, "A linear algorithm for optimization over directed graphs with geometric convergence," *IEEE Control Syst. Lett.*, vol. 2, no. 3, pp. 315–320, 2018.

[27] A. Nedić, A. Olshevsky, and W. Shi, "Achieving geometric convergence for distributed optimization over time-varying graphs," *SIAM J. Optim.*, vol. 27, no. 4, pp. 2597–2633, 2017.

[28] B. Gharesifard, and J. Cortés, "Distributed strategies for generating weight-balanced and doubly stochastic digraphs," *Eur. J. Control*, vol. 18, no. 6, pp. 539–557, 2012.

[29] A. I. Rikos, T. Charalambous, and C. N. Hadjicostis, "Distributed weight balancing over digraphs," *IEEE Trans. Control Netw. Syst.*, vol. 1, no. 2, pp. 190–201, 2014.

[30] S. Pu, "A robust gradient tracking method for distributed optimization over directed networks," in *IEEE Conf. Decis. Control (CDC)*, 2020, pp. 2335–2341.

[31] J. Xu, S. Zhu, Y. C. Soh, and L. Xie, "Convergence of asynchronous distributed gradient methods over stochastic networks," *IEEE Trans. Autom. Control*, vol. 63, no. 2, pp. 434–448, 2018.

[32] Y. Kajiyama, N. Hayashi, and S. Takai, "Linear convergence of consensus-based quantized optimization for smooth and strongly convex cost functions," *IEEE Trans. Autom. Control*, vol. 66, no. 3, pp. 1254–1261, 2021.

[33] R. A. Horn, and C. R. Johnson, "Matrix analysis," *Cambridge university press*, 2012.

[34] G. Qu, and N. Li, "Harnessing smoothness to accelerate distributed optimization," *IEEE Trans. Control Netw. Syst.*, vol. 5, no. 3, pp. 1245–1260, 2018.